Learning Feasibility to Imitate Demonstrators with Different Dynamics

Zhangjie Cao¹, Yilun Hao¹, Mengxi Li², Dorsa Sadigh^{1,2}

¹Department of Computer Sciences, Stanford University, United States

²Department of Electrical Engineering, Stanford University, United States
caozj@cs.stanford.edu, {yilunhao,mengxili}@stanford.edu, dorsa@cs.stanford.edu

Abstract: The goal of learning from demonstrations is to learn a policy for an agent (imitator) by mimicking the behavior in the demonstrations. Prior works on learning from demonstrations assume that the demonstrations are collected by a demonstrator that has the same dynamics as the imitator. However, in many real-world applications, this assumption is limiting — to improve the problem of lack of data in robotics, we would like to be able to leverage demonstrations collected from agents with different dynamics. This can be challenging as the demonstrations might not even be feasible for the imitator. Our insight is that we can learn a feasibility metric that captures the likelihood of a demonstration being feasible by the imitator. We develop a feasibility MDP (f-MDP) and derive the feasibility score by learning an optimal policy in the f-MDP. Our proposed feasibility measure encourages the imitator to learn from more informative demonstrations, and disregard the far from feasible demonstrations. Our experiments on four simulated environments and on a real robot show that the policy learned with our approach achieves a higher expected return than prior works. We show the videos of the real robot arm experiments on our website.

Keywords: Imitation Learning, Learning from Agents with Different Dynamics

1 Introduction

Imitation learning aims to learn a well-performing policy from demonstrations. Standard imitation learning algorithms usually assume that the *demonstrator* (the agent that generates the demonstrations) and the imitator (the agent that is learning a policy) share the same dynamics, i.e., the transition functions are the same [1, 2, 3, 4]. Specifically, in a given state, with the same action, both the demonstrator and the imitator transition to the same distribution of next states. However, this assumption limits the usage of already collected demonstrations. Imagine a setting, where a set of demonstrations are collected on a 7 Degrees of Freedom (DoF) robot arm shown in Fig. 1 to place a book on the empty area of the shelf (on the left) without colliding with the books that are already placed on the right side of the shelf. Later, we might decide to buy a different arm with 3 DoF (e.g., only the joints circled in green as shown in the figure are used). We would like to learn a policy for this 3 DoF robot arm that can achieve the same task—placing the book on the empty region of the shelf—using the originally collected demonstrations on the 7 DoF arm. In general, our goal is to enable using and reusing data collected on robots with different dynamics or embodiments to tackle the problem of lack of in-domain data in robotics. The 3 DoF robot arm should still be able to learn a policy based on feasible or nearly feasible demonstrations from an agent with different dynamics, e.g., using the trajectories that go over the bookshelf in Fig. 1. Motivated by this example, we relax the assumption of shared dynamics between the imitator and demonstrator so that the data can be collected from demonstrators with the same state space but different dynamics from the imitator, e.g., demonstrators with different embodiments, body schemas, joints, or rigid body structures.

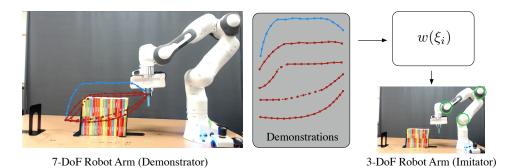


Figure 1: An example of imitating demonstrators with feasibility. The left image shows that a set of demonstrations (blue and red trajectories) are available for the 7 DoF robot arm. We aim to learn a policy for the 3 DoF robot (joints are circled in green) by learning from the demonstrations of the 7 DoF robot (blue is feasible and red is infeasible). We learn a feasibility score to reweight each demonstration to conduct imitation learning.

Prior works in imitation learning from demonstrators with different dynamics typically rely on state-only demonstrations and learn a policy to maximally follow the sequence of states in demonstrations [5, 6]. Such learning techniques assume that all of the collected demonstrations are useful for the imitator. However, it is possible that demonstrations drawn from agents with other dynamics can be useless or even harmful for the imitator because they may not be feasible for the imitator. Going back to the example in Fig. 1, the red trajectories that move around the stack of books are not feasible for the 3 DoF robot arm. Imitating such trajectories may cause the 3 DoF robot arm to maximally follow these trajectories and even collide with the existing stack of books. Therefore, it is crucial to identify and avoid trajectories that are far from feasible for the imitator, and instead learn more from useful demonstrations, e.g., the blue trajectories that go over the shelf that are still feasible for a robot with 3 DoF.

To avoid the influence of useless or harmful demonstrations from agents with different dynamics, we rely on a feasibility score, which measures how feasible a trajectory is for the imitator, and select trajectories with high feasibility to imitate. For example, the blue trajectories should have higher feasibility than the red trajectories in Fig. 1. Prior work such as Cao and Sadigh [7] estimate the feasibility score by computing the distances of demonstrations and corresponding trajectories but the performance highly relies on the accuracy of the inverse dynamics model, which can be difficult to learn. Our key idea is to directly learn a feasibility score for the imitator based on the collected demonstrations. Specifically, we model the imitator environment as an MDP and build a feasibility Markov Decision Process (f-MDP) based on the imitator's MDP and the trajectories provided by the demonstrator. The optimal policy for the f-MDP maximally follows the behavior of the demonstrations but is limited by the imitator's environment. This optimal policy helps assign a feasibility score over the demonstrations. We conduct imitation learning on the demonstrations re-weighted by the feasibility score to learn the final policy for the imitator. We experiment with several simulation environments and a manipulation task with a Panda Franka arm. We show that the policy learned from demonstrations re-weighted by our feasibility achieves higher performance compared to other methods.

2 Related Works

Imitation Learning. Imitation learning seeks a policy that best imitates demonstrations. Current imitation learning methods can be roughly divided into Behavior Cloning (BC), Inverse Reinforcement Learning (IRL) and Generative Adversarial Imitation Learning (GAIL). BC directly learns the policy from a sequence of state-action pairs via supervised learning [8], where dataset aggregation [9] or policy aggregation [10, 11] are proposed to address the compounding errors problem. IRL first learns a reward function that best matches demonstrations and then finds a policy through reinforcement learning to maximize the recovered reward [1, 12, 2, 13]. GAIL learns the expert policy by matching the occupancy measure between the policy and the demonstrations [4].

However, most imitation learning works require that the demonstrations consist of a sequence of states and actions. When only state observations are available, new imitation learning algorithms are proposed to address the lack of actions. Torabi et al. [14] recover the actions between consecutive states through an inverse dynamics model. GAIL-based works directly match the state occupancy measure between the demonstrations and the policy [15, 16, 17]. However, imitation learning methods learned from either state-action or state-only demonstrations assume that the demonstrator and the imitator have the same dynamics. Since demonstrations from different dynamics may not be feasible for the imitator, directly imitating cannot achieve the same optimal behavior, and may cause unknown suboptimal outcomes. Thus, standard imitation learning algorithms do not fit our problem setting.

Learning from Demonstrations with Different Dynamics. Early works model this problem as a correspondence problem between the demonstrator and the imitator, and map states and actions in demonstrations to the imitator's states and actions [18, 19]. Englert et al. [20] align the state trajectory distributions to address the correspondence problem. Calinon et al. [21] model the demonstrations as a Gaussian mixture model within a projected lower-dimensional subspace. Eppner et al. [22] learn a task description. Domain randomization methods learn the correspondence as an invariant latent space by randomizing domains [23, 24, 25]. Zhang et al. [26] learns a translation mapping to model the correspondence. However, modeling correspondence requires that there exists a strict correspondence between the MDP of the demonstrator and the imitator. Recent works instead only assume the shared state space between the demonstrator and the imitator, and address the different dynamics problem by encouraging the imitator to maximally follow the state trajectory of the demonstrator [5, 27, 28, 6]. However, all these works ignore an important challenge—that is the demonstrations may be far from feasible for the imitator. Enforcing the imitator to follow such trajectories may lead to unknown behavior. We focus on this challenge and develop a feasibility score to down-weight demonstration trajectories that are far from feasible for the imitator. Compared to the works that learn feasibility to filter infeasible trajectories [7], we do not require the inverse dynamics model, which can make our setting more generalizable to different environments.

3 Problem Statement

In our problem setting, an imitator aims to learn from demonstrations collected from N demonstrators with various dynamics. We formalize the demonstrators and the imitator each as a standard Markov decision process (MDP). For each demonstrator j, $(1 \leq j \leq N)$, the MDP is formalized as $\mathcal{M}_j^d = \langle \mathcal{S}, \mathcal{A}_j^d, p_j^d, \mathcal{R}, \rho_0, \gamma \rangle$. The MDP for the imitator is $\mathcal{M}^i = \langle \mathcal{S}, \mathcal{A}^i, p^i, \mathcal{R}, \rho_0, \gamma \rangle$. \mathcal{S} is the shared state space for all environments. \mathcal{A}_j^d and \mathcal{A}^i are the action spaces and $p_j^d: \mathcal{S} \times \mathcal{A}_j^d \times \mathcal{S} \to [0,1]$ and $p^i: \mathcal{S} \times \mathcal{A}^i \times \mathcal{S} \to [0,1]$ are the transition probabilities for each demonstrator and the imitator respectively. Note that in our problem setting, we use the transition function p to denote dynamics and the demonstrators and the imitator may have different dynamics and action spaces. ρ_0 is the shared initial state distribution for all MDPs. $\mathcal{R}: \mathcal{S} \times \mathcal{S} \to \mathbb{R}$ is the reward function. Note that we make the assumption that the reward function is based on state transitions and is shared between the demonstrators and the imitator, which is a common assumption used in prior work [5, 7], and is usually satisfied since the demonstrators and the imitator conduct the same task in the same context. γ is the shared discount factor. A policy for the imitator $\pi^i: \mathcal{S} \times \mathcal{A}^i \to [0,1]$ defines a probability distribution over the space of actions in a given state. An optimal policy π^* maximizes the expected return $\eta_{\pi^i} = \mathbb{E}_{s_0 \sim \rho_0, \pi^i} [\sum_{t=0}^{\infty} \gamma^t \mathcal{R}(s_t, a_t, s_{t+1})]$, where t indicates the time step.

We aim to learn a policy π^i for the imitator, given a set of demonstrations from different demonstrators $\Xi^j = \{\xi_1^j, \dots, \xi_D^j\}_{j \in \{1...N\}}$ where each trajectory is a sequence of states $\xi = \{s_0^d, s_1^d, \dots, s_H^d\}$. We assume that the optimal policy can be learned by imitating the useful demonstrations, which is a general assumption adopted by prior imitation learning works [8, 4, 5, 7]. The violation of this assumption, as shown in prior works, leads to learning a suboptimal policy. Note that we discard actions from the demonstrations instead of imitating the state-action trajectories because different action spaces between the demonstrators and the imitator make it impossible to imitate the actions.

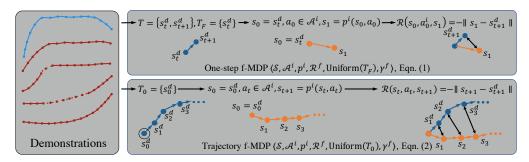


Figure 2: The illustration and comparison of one-step f-MDP and trajectory f-MDP. The blue state transition and trajectory are from the demonstrations while the orange state transition and trajectory are rollouts in the f-MDP. One-step f-MDP collects the states in the *Former Set* and uses the uniform distribution over the states as the initial state distribution. Trajectory f-MDP collects the initial state of all the demonstrations and uses a uniform distribution over them as the initial state distribution.

Challenges. The core challenges of imitation learning from demonstrations with different dynamics are: (1) How to imitate useful demonstrations with different dynamics, (2) How to avoid harmful demonstrations misleading the imitator. Prior works have studied and made progress for the first challenge [5, 6], but the second challenge is still under-studied. Strong assumptions such as access to or learning an accurate inverse dynamics model are needed to filter out harmful demonstrations [7]. We address the second challenge by learning a feasibility score that measures how likely it is for a demonstration to be feasible for the imitator with minimal assumptions: only using the environment of the imitator, i.e., we can collect interaction data in the environment but we do not know the exact reward and transition function of the imitator.

4 Feasibility-Based Imitation Learning

The feasibility of a trajectory depends on the feasibility of each state transition in the trajectory, i.e., if (s_t, s_{t+1}) is feasible for all time steps. A state transition (s_t, s_{t+1}) is feasible when there exists an action $a_t^i \in A^i$ such that $p^i(s_t, a_t^i, s_{t+1}) = 1$ for deterministic transitions or $p^i(s_t, a_t^i, s_{t+1}) > 0$ for stochastic transitions. In this section, we discuss the deterministic MDP setting and discuss the stochastic setting in Appendix.

Feasibility can be directly measured by a perfect inverse dynamics model $f: \mathcal{S} \times \mathcal{S} \to \mathcal{A}$ that takes a state transition $(s_t, s_{t+1}) \in \mathcal{S} \times \mathcal{S}$ as the input and outputs the action $a_t \in \mathcal{A}$ that achieves the transition if feasible or outputs 'Infeasible'. However, having access to this model is often non-trivial and such a binary feasibility measurement as f discards all infeasible demonstrations without considering any useful information from slightly infeasible trajectories.

Our goal is to learn a policy $\pi^i:\mathcal{S}\to\mathcal{A}^i$ for the imitator to maximally achieve the state transitions in the demonstrations. This means that if the state transition (s^d_t, s^d_{t+1}) from a demonstration is feasible, the next state produced by π^i , i.e., $s^i_{t+1}=p^i(s^d_t,\pi^i(s^d_t))$ should be equal to s^d_{t+1} . Otherwise, we would like the policy to output an action that ensures the next state s^i_{t+1} is as close as possible to the next state from the demonstration s^d_{t+1} . Therefore, the distance between s^i_{t+1} and s^d_{t+1} can serve as a measure of feasibility, where a smaller distance corresponds to a higher likelihood of feasibility. To learn this policy, we design a feasibility MDP (f-MDP), where we ensure that the optimal policy of the f-MDP satisfies the above requirement. f-MDP is defined as $M^f = \langle \mathcal{S}, \mathcal{A}^i, p^i, \mathcal{R}^f, \rho^f_0, \gamma^f \rangle$. We will now discuss our choices for the components of f-MDP.

One-step f-MDP. First, recall that our goal is to learn a policy for the imitator to maximally achieve the state transitions in the demonstrations. So the policy should be learned in an environment with the same state-action space and transition probability as the imitator. We would like the reward of the f-MDP to encourage maximally achieving the state transitions in the demonstrations. Let us first collect all the state transitions $T = \{(s_t^d, s_{t+1}^d)\}$ in all of the demonstrations. We define the Former Set to be the set of states in the demonstrations that one can transition from: $T_F = \{s_t^d : (s_t^d, s_{t+1}^d) \in T\}$. The

initial state distribution ρ_0^f can be defined uniformly over the Former Set as Uniform (T_F) . Here, we assume that all the states in the Former Set can be visited by the imitator. We define the reward of a *One-step f-MDP* so that it matches the one-step transitions from the Former Set:

$$s_t^d \sim \text{Uniform}(T_F), \quad s = s_t^d, \quad s' = p^i(s, a), \quad R^f(s, a, s') = -f_{\text{dis}}(s', s_{t+1}^d),$$
 (1)

where (s_t^d, s_{t+1}^d) is a state transition in the demonstrations and $a \in A^i$ is sampled from the action space of the f-MDP. f_{dis} is a function that measures the distance between the states (e.g., the L2 distance). We define the reward to penalize the distance between s' and s_{t+1}^d .

Trajectory f-MDP. The one-step f-MDP suffers from an important shortcoming: the assumption that all states in the Former Set must be visited by the imitator can be violated, because the demonstrators have different dynamics from the imitator and some demonstration states can never be reached by the imitator. So we cannot set $Uniform(T_F)$ as the initial state distribution for the f-MDP. We instead collect the initial state s_0^d of all the demonstrations, $T_0 = \{s_0^d\}$, and define the initial state distribution of the $Trajectory\ f\text{-}MDP$ as $Uniform(T_0)$. Since all the demonstrators and the imitator share the initial state distribution, all states in T_0 can be visited by the imitator. We define the reward as:

$$s_0^d \sim \text{Uniform}(T_0), \quad s_0 = s_0^d, \quad s_{t+1} = p^i(s_t, a), \quad R^f(s_t, a, s_{t+1}) = -f_{\text{dis}}(s_{t+1}, s_{t+1}^d), \quad (2)$$

We use the L2 distance for f_{dis} in our experiments. Similar to the one-step f-MDP $a \in A^i$ is sampled from the action space of the imitator.

Learning Feasibility. Given the Trajectory f-MDP defined above, for each demonstration trajectory ξ , the highest reward achieved in this f-MDP reflects the feasibility score of the trajectory. We use reinforcement learning to learn the optimal policy of the Trajectory f-MDP, π^* . We then derive the feasibility of each demonstration trajectory ξ as a function of the trajectory f-MDP reward:

$$w(\xi) = \exp\left(\frac{-\sum_{t=1}^{N} (\gamma^f)^t f_{\text{dis}}(s_t, s_t^d) - C}{\sigma}\right).$$
 (3)

 s_t is the state at step t in the rollout derived by the policy π^* . We use an exponential function of the cumulative reward since the cumulative reward is always negative and the exponential function can bound the feasibility in the range of [0,1]. The parameter C is used to shift the function to avoid the situation where the cumulative reward is extremely negative, while the parameter σ controls how low the reward can be, and when a demonstration can be fully filtered out by assigning a feasibility of close to 0. In practice, C is usually set as the maximal cumulative reward over all demonstrations to ensure the maximal feasibility is 1.

For the feasibility of each state transition (s_t^d, s_{t+1}^d) , we use the feasibility of the trajectory it belongs to: $w((s_t^d, s_{t+1}^d)) = w(\xi_i)$, where $(s_t^d, s_{t+1}^d) \in \xi_i$. We do not use the state distance at each time step between s_{t+1} and s_{t+1}^d as in the One-step f-MDP because such measurement suffers from the fact that within a trajectory, the reward of later steps are influenced by former steps. For example, if s_t diverges from s_t^d , s_{t+1} will diverge more from s_{t+1}^d . So the per-step reward is an unfair measure of feasibility for the state transition (s_t^d, s_{t+1}^d) at different time steps t. Therefore, we use the accumulative reward of the whole trajectory as our feasibility measure, where all the state transitions share the same value.

The discount factor γ^f is usually set as $\gamma^f < 1$ to reduce compounding errors. Specifically, the length of a rollout in the f-MDP is the same as the corresponding demonstration, which can be very long. If the state in the rollout starts to diverge from the demonstration trajectory at t, meaning that $\|s_t - s_t^d\| > 0$, the steps after time step t even diverge more from the demonstration. This makes the trajectory reward for all the infeasible trajectories very low and does not allow for discriminating among different infeasible trajectories. Therefore, we set a discount factor of $\gamma^f < 1$ to discount or even ignore the trajectory reward at later steps.

Leveraging our Trajectory f-MDP design, feasible trajectories still receive the maximal reward of 0 since each state in the rollout will perfectly match the demonstration thus having a feasibility of 1 as in Eqn. (3). Instead, infeasible trajectories receive negative rewards leading to smaller feasibility scores, which reflects how far away the demonstration is from the closest feasible trajectory.

One may worry about the time complexity of our approach since we need an additional RL training to learn an optimal policy for the f-MDP. However, the f-MDP is a lot simpler compared to the imitator's MDP since the initial distribution is reduced from the distribution of all possible states in the demonstration set to a discrete distribution over the initial states of the demonstrations. This can simplify the time complexity of finding the optimal policy for the f-MDPs.

Algorithm. Using the feasibility metric in Eqn. (3), we assign each state transition with the same feasibility of the trajectory it belongs to. Directly weighing the imitation loss as [29] may lead to gradients that are close to 0 if a batch of data all have low feasibility. This can make the algorithm inefficient by wasting samples from many iterations. Instead, for a more efficient training, we define a discrete probability distribution p_w over the collection of state transitions in all the demonstrations: T, where the probability of a state-transition (s_t^d, s_{t+1}^d) as $p_w((s_t^d, s_{t+1}^d)) = \frac{w((s_t^d, s_{t+1}^d))}{\sum_{s_t^d, s_{t+1}^d} p_{s_t^d}} \frac{w((s_t^d, s_{t+1}^d))}{\sum_{s_t^d, s_t^d} p_{s_t^d}} \frac{w((s_t^d, s_t^d))}{\sum_{s_t^d, s_t^d} p_{s_t^d}} \frac{w((s_t^d, s_t^d))}{\sum_{s_t^d} p_{s_t^d}} \frac{w((s_t^d, s_t^d))}{\sum_{s_t^d} p_{s_t^d}} \frac{w((s_t^d, s_t^d))}{\sum_{s$

State transitions with larger feasibility will be sampled more often. Using the sampling distribution p_w , we can embed our method into any imitation learning algorithm to enable learning from demonstrations with different dynamics. We show the algorithm block in the Appendix.

Sampling More Demonstrations with the Feasibility Score. When the existing useful demonstrations are too scarce to learn a well-performing imitation learning policy, we need to acquire more demonstrations from the demonstrators. But collecting new demonstrations can be expensive, so we often can only acquire a limited budget of demonstrations. We thus need to collect the most useful demonstrations within this limited budget. The proposed feasibility metric provides a criterion to decide the similarity between the imitator and each demonstrator. If a demonstrator has a higher similarity, we sample more from this demonstrator because its demonstrations are more likely to be feasible. Specifically, we create a probability distribution p_i over all demonstrators:

$$p_{j} = \frac{\frac{1}{|\Xi^{j}|} \sum_{\xi^{j} \in \Xi^{j}} w(\xi^{j})}{\sum_{j=1}^{N} \frac{1}{|\Xi^{j}|} \sum_{\xi^{j} \in \Xi^{j}} w(\xi^{j})}.$$
 (4)

We repeatedly and independently query the demonstrator j according to p_j and collect a demonstration. The proposed sampling strategy samples more demonstrations from closer demonstrators. We empirically show that the sampling strategy derived from our feasibility performs better than uniform sampling or using other feasibility metrics as in [7].

5 Experimental Results

We experiment with four MuJoCo environments, a simulated Franka Panda Arm, and a real Franka Panda Arm. We also show results on various compositions of demonstrations of different dynamics and the performance gain when we are given a larger budget to collect demonstrations. We compare our approach with a standard imitation learning algorithm: GAIL [4], imitation learning from demonstrations with different dynamics methods without a measure of feasibility: SAIL [5], and with a feasibility score: ID-Feas [7], which uses an inverse dynamics model to estimate feasibility.

5.1 MuJoCo Experiments

Swimmer. The swimmer agent has three links and two joints. The goal of the agent is to move forward by rotating the joints. As shown in Fig. 3(a), we create different dynamics by setting the joint limit of the front and the back joints, denoted by (α_f, α_b) . The original Swimmer environment has $(\alpha_f, \alpha_b) = (100^\circ, 100^\circ)$. We create four demonstrator environments (α_f, α_b) : (i) $(100^\circ, 12^\circ)$, (ii) $(100^\circ, 20^\circ)$, (iii) $(100^\circ, 100^\circ)$, and (iv) $(10^\circ, 100^\circ)$. We also create the imitator environment by setting $(\alpha_f, \alpha_b) = (100^\circ, 10^\circ)$. The demonstrators (i) and (ii) are closer to the imitator in terms of their dynamics, while the demonstrators (iii) and (iv) are farther.

Walker2d. The Walker2d is an agent with two legs where each leg consists of 3 joints. We create different dynamics by using different frictions β for the feet, i.e., the link that touches the ground. The original Walker2d uses $\beta = 0.9$. We create two settings to show high friction and low friction of the

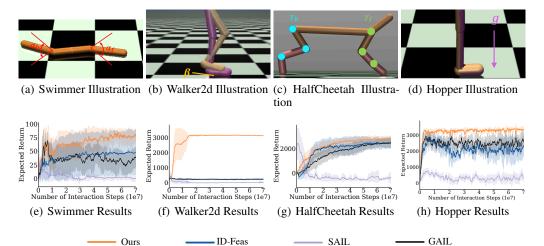


Figure 3: Illustration of different dynamics in (a) Swimmer: varying the joint limit of the front and back joints $(\alpha_f \text{ and } \alpha_b)$. (b) Walker2d: varying the friction of the feet (β) . (c) HalfCheetah: varying the joint control force of the front and back joints by multiplying a factor γ_f and γ_b with the front and back joint force. (d) Hopper: varying the gravitational constant respectively (e-h) show the expected return in these four environments.

imitator with a mix of frictions for the demonstrators. In the first setting, there are four demonstrators: (i) $\beta=19.9$, (ii) $\beta=9.9$, (iii) $\beta=0.9$, and (iv) $\beta=0.7$. The imitator has $\beta=24.9$. In the second setting, there are four demonstrators: (i) $\beta=29.9$, (ii) $\beta=19.9$, (iii) $\beta=1.1$, and (iv) $\beta=0.7$. The imitator has $\beta=0.9$.

HalfCheetah. The HalfCheetah is an agent with two legs at the front and back of the body, where each leg consists of three joints. We create different dynamics by varying the control force limit of joints of the front leg and back leg, where we multiply a factor γ_f with the original control force limit of the front leg and multiply γ_b with the limit of the back leg. We create two settings, where the demonstrators have low and high similarity with each other. In the first setting, there are four demonstrators with (γ_f, γ_b) : (i) (0.05, 1), (ii) (0.5, 1), (iii) (1, 0.5), and (iv) (1, 0.05). The imitator has $(\gamma_f, \gamma_b) = (0.01, 1)$. In the second setting, there are four demonstrators with (γ_f, γ_b) : (i) (0.01, 1), (ii) (0.05, 1), (iii) (1, 0.05), and (iv) (1, 0.01). The imitator has $(\gamma_f, \gamma_b) = (0.01, 1)$.

Hopper. The Hopper is an agent with one leg consisting of 3 joints. We create different dynamics by using different gravitational constants g. The original Hopper uses g=9.81. We create four demonstrator environments: (i) g=20.0, (ii) g=9.81, (iii) g=5.0, and (iv) g=2.0. We also create the imitator environment by setting g=15.0.

The detailed composition of demonstrations for all four environments is included in the Appendix. For all the Mujoco environments, we evaluate the expected return of each policy by rolling out 100 trajectories in the environment with the policy and compute the average expected return of the 100 trajectories. We run each experiment for 5 times and show the mean and the standard deviation.

Results. We show the expected return with respect to the number of steps for the four different environments in Fig. 3. We show the results of the second setting for the Walker2d and the HalfCheetah in the Appendix. We observe that our proposed feasibility achieves the best performance among all the methods. The highest p-value comparing our method to baselines is 0.116 with ID-Feas for Swimmer, 2.55e-14 with GAIL for Walker2d, 0.188 with ID-Feas for HalfCheetah, and 0.026 with GAIL for Hopper. In particular, our method outperforms ID-Feas, which indicates that the proposed feasibility more accurately filters out far from feasible demonstrations. SAIL performs even worse than GAIL, this is because SAIL can more strictly follow the state sequences of demonstrations than GAIL including those far from feasible demonstrations. Our demonstration set is composed of a high percentage of demonstrations from unrelated dynamics, which can mislead SAIL's learned policy.

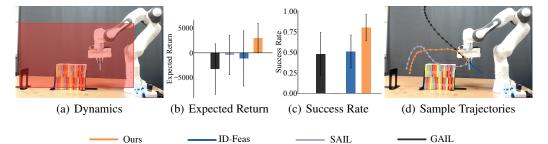


Figure 4: (a) Illustration of different dynamics in the real robot arm environment. (b-c) The bar plots show the expected return and success rate compared to other methods. (d) Sampled trajectories using different methods.

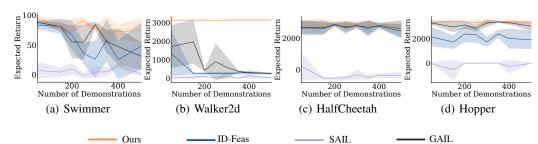


Figure 5: (a-d) The expected return when increasing the number of demonstrations from agents with unrelated dynamics. The results in Fig. 3 correspond to using 500 demonstrations from each unrelated dynamics. In both of these settings, there will also be a fixed number of demonstrations from agents with related dynamics as shown in Appendix.

5.2 Simulated and Real Robot Arm Experiments

Setup. We create a simulated robot arm based on a Panda Robot Arm implemented in the PyBullet [30] and a real robot arm environment using a Franka Panda Arm¹. We include the results for the simulated robot arm in the Appendix. As shown in Fig. 4(a), we create a task of moving a book to the shelf but the closest region on the shelf is full. So we need to move the book to the empty area of the shelf without colliding with the shelf and the existing books on the shelf. We create two dynamics for the robot arm: using a 7-DoF control which can move freely in the 3D space, and using a 3-DoF control, which is limited to moving on the red plane area. We collect demonstrations from both 7-DoF and 3-DoF controllers and aim to learn an optimal policy for the 3-DoF robot.

For evaluation, we use two metrics: (1) The expected return based on a reward penalizing collision with the shelf and existing books while rewarding the successfully moving the book to the empty area of the shelf within the time limit. More detail on the reward is in the Appendix. (2) The success rate of finishing the task over 100 trials.

We observe that the proposed approach outperforms the baseline methods both in expected return and success rate as shown in Fig. 4. The highest p-value for the expected return is 2.432×10^{-7} and for the success rate is 3.534×10^{-8} (both with ID-Feas), which are statistically significant.

5.3 Analysis

We conduct experiments with varying compositions of demonstrations and investigate the performance of different approaches when we have the budget to acquire additional demonstrations. We show the results of varying the number of demonstrations from all demonstrators in the Appendix.

Varying the Number of Demonstrations from each Unrelated Demonstrator. For the first three experiment settings in the Mujoco environment, we have two demonstrators with similar dynamics to the imitator and two demonstrators with far apart dynamics. We vary the number of demonstrations

¹https://www.franka.de

from the far apart demonstrators to investigate their influence on the different methods. We conduct experiments on the first setting for the Swimmer, Walker2D, HalfCheetah, and Hooper and report the results in Fig. 5(a), 5(b), 5(c) and 5(d). With an increasing number of demonstrations from the far apart demonstrators, the expected return of all the methods decreases, while our method shows the best performance consistently across different numbers of demonstrations. This demonstrates that our feasibility can effectively filter out far from feasible demonstrations and ensure the policy learns from useful demonstrations.

6 Conclusion

Summary. We propose an algorithm to learn a feasibility metric to imitate demonstrations drawn from agents with different dynamics. Our feasibility metric captures how likely it is for each demonstration to be feasible for the imitator. We develop a feasibility MDP (f-MDP) and derive the feasibility by learning the optimal policy for the f-MDP. We show that the policy learned from the demonstrations reweighted by the proposed feasibility score outperforms other imitation learning methods in various environments and different compositions of demonstrations.

Limitations and Future Work. Our work only addresses the problem of filtering out far from feasible demonstrations, but does not solve the problem of learning a policy from those feasible but suboptimal or nearly feasible demonstrations from different dynamics. There are situations where demonstrations are feasible but not optimal for the imitator, especially when the ability of the demonstrator is more restricted than the imitator. In the future, we aim to study these more general settings.

Acknowledgments

We would like to thank FLI grant RFP2-000, NSF Awards 1941722 and 1849952, and DARPA HiCoN project for their support of this project.

References

- [1] P. Abbeel and A. Y. Ng. Apprenticeship learning via inverse reinforcement learning. In *ICML*, pages 1–8, 2004.
- [2] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey. Maximum entropy inverse reinforcement learning. In *AAAI*, 2008.
- [3] A. H. Qureshi, B. Boots, and M. C. Yip. Adversarial imitation via variational inverse reinforcement learning. In *ICLR*, 2019.
- [4] J. Ho and S. Ermon. Generative adversarial imitation learning. In *NeurIPS*, volume 29, 2016.
- [5] F. Liu, Z. Ling, T. Mu, and H. Su. State alignment-based imitation learning. In *ICLR*, 2019.
- [6] I. Radosavovic, X. Wang, L. Pinto, and J. Malik. State-only imitation learning for dexterous manipulation. *arXiv preprint arXiv:2004.04650*, 2020.
- [7] Z. Cao and D. Sadigh. Learning from imperfect demonstrations from agents with varying dynamics. *IEEE Robotics and Automation Letters (RA-L)*, 2021.
- [8] M. Bain and C. Sammut. A framework for behavioural cloning. In *Machine Intelligence 15*, 1995.
- [9] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *AISTATS*, 2011.
- [10] H. Daumé, J. Langford, and D. Marcu. Search-based structured prediction. *Machine learning*, 2009.
- [11] S. Ross and D. Bagnell. Efficient reductions for imitation learning. In AISTATS, 2010.
- [12] A. Y. Ng, S. J. Russell, et al. Algorithms for inverse reinforcement learning. In *ICML*, pages 663–670, 2000.
- [13] C. Finn, S. Levine, and P. Abbeel. Guided cost learning: Deep inverse optimal control via policy optimization. In *International conference on machine learning*, pages 49–58. PMLR, 2016.
- [14] F. Torabi, G. Warnell, and P. Stone. Behavioral cloning from observation. In *IJCAI*, pages 4950–4957, 2018.
- [15] Y. Schroecker and C. L. Isbell. State aware imitation learning. In *NeurIPS*, 2017.
- [16] F. Torabi, G. Warnell, and P. Stone. Generative adversarial imitation from observation. *Imitation, Intent, and Interaction (13) Workshop at ICML*, 2019.
- [17] W. Sun, A. Vemula, B. Boots, and D. Bagnell. Provably efficient imitation learning from observation alone. In *ICML*, 2019.
- [18] C. L. Nehaniv, K. Dautenhahn, et al. The correspondence problem. *Imitation in animals and artifacts*, 41, 2002.
- [19] B. D. Argall, S. Chernova, M. Veloso, and B. Browning. A survey of robot learning from demonstration. *Robotics and autonomous systems*, 2009.
- [20] P. Englert, A. Paraschos, J. Peters, and M. P. Deisenroth. Addressing the correspondence problem by model-based imitation learning. In *ICRA Workshop on Autonomous Learning*, 2013.

- [21] S. Calinon, F. Guenter, and A. Billard. On learning, representing, and generalizing a task in a humanoid robot. TSMC, 2007.
- [22] C. Eppner, J. Sturm, M. Bennewitz, C. Stachniss, and W. Burgard. Imitation learning with generalized task descriptions. In *ICRA*, 2009.
- [23] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In 2017 IEEE/RSJ international conference on intelligent robots and systems (IROS), pages 23–30. IEEE, 2017.
- [24] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In 2018 IEEE international conference on robotics and automation (ICRA), pages 3803–3810. IEEE, 2018.
- [25] S. Zakharov, W. Kehl, and S. Ilic. Deceptionnet: Network-driven domain randomization. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 532–541, 2019.
- [26] Q. Zhang, T. Xiao, A. A. Efros, L. Pinto, and X. Wang. Learning cross-domain correspondence for control with dynamics cycle-consistency. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=QIRlze316hX.
- [27] P. Sharma, L. Mohan, L. Pinto, and A. Gupta. Multiple interactions made easy (mime): Large scale demonstrations data for imitation. In *Conference on robot learning*, pages 906–915. PMLR, 2018.
- [28] T. Gangwani and J. Peng. State-only imitation with transition dynamics mismatch. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=HJgLLyrYwB.
- [29] Y.-H. Wu, N. Charoenphakdee, H. Bao, V. Tangkaratt, and M. Sugiyama. Imitation learning from imperfect demonstration. In *ICML*, 2019.
- [30] E. Coumans and Y. Bai. Pybullet, a python module for physics simulation for games, robotics and machine learning. 2016.
- [31] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *ICML*, 2015.
- [32] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016.

In this Appendix, we provide the proposed method under stochastic MDPs in Sec. A. We provide the details of the algorithm in Sec. B. We further explain our experimental details in Sec. C. Finally, we show some extra experimental results in Sec. D.

A The Extension to Stochastic MDPs

In our main text, we discussed the deterministic MDP setting as all our experiments are in the deterministic setting, but here we would like to extend the discussion to the stochastic MDP case. In a stochastic MDP, our goal is still to learn a policy $\pi^i:\mathcal{S}\to\mathcal{A}^i$ for the imitator to maximally achieve the state transitions in the demonstrations. This means that if the state transition (s^d_t,s^d_{t+1}) from a demonstration is more likely to be feasible, the expected distance between the next state s^i_{t+1} produced by π^i and s^d_{t+1} should be small, i.e., $\mathbb{E}_{s^i_{t+1}\sim p^i(s^d_t,\pi^i(s^d_t))}[f_{\mathrm{dis}}(s^i_{t+1},s^d_{t+1})]$ should be small. Therefore, the expected distance between s^i_{t+1} and s^d_{t+1} can serve as a measure of feasibility, where a smaller distance corresponds to a higher feasibility.

Under a feasibility metric defined by the expected distance, we can use the same design of f-MDP as the deterministic MDP case: $M^f = \langle \mathcal{S}, \mathcal{A}^i, p^i, \mathcal{R}^f, \rho_0^f, \gamma^f \rangle$. The state space, the action space and the transition probability are all the same as the imitator's. The initial state distribution is defined as Uniform (T_0) . The reward is defined as:

$$s_0^d \sim \text{Uniform}(T_0), \quad s_0 = s_0^d, \quad s_{t+1} = p^i(s_t, a), \quad R^f(s_t, a, s_{t+1}) = -f_{\text{dis}}(s_{t+1}, s_{t+1}^d).$$
 (5)

Maximizing the expected return in the f-MDP in such a design will minimize the expected state distance between the learned policy and the demonstrations, which matches our definition of feasibility. After we learn the optimal policy π^* for the f-MDP, we can derive the feasibility for each trajectory ξ with the expected state distance between the demonstration and the policy:

$$w(\xi) = \exp\left(\frac{-\mathbb{E}_{s_t \sim \pi^*} \left[\sum_{t=1}^N (\gamma^f)^t f_{\text{dis}}(s_t, s_t^d)\right] - C}{\sigma}\right).$$
 (6)

In our experiments, we only consider the case where the MDP is deterministic.

B Algorithm

We go through the steps of the algorithm of learning feasibility to imitate demonstrators with different dynamics in Algorithm 1. For the state-based imitation learning algorithm used to learn the final policy from reweighted demonstrations, we use the state-based GAIL as [7]. Lines 1-6 show the process of constructing N f-MDPs (one for each demonstrator) and training the optimal policy for each f-MDP. Line 7 shows how to compute feasibility with the optimal policy for each f-MDP. Lines 8-11 show imitation learning over the newly reweighted demonstrations.

Algorithm 1: Algorithm

Input: Demonstrations Ξ^j from each demonstrator \mathcal{M}_j^d , $(1 \le j \le N)$.

Construct the trajectory f-MDP M_i^f based on the demonstration set Ξ_j according to Eqn. (2);

Train an optimal policy π_j^* for the trajectory f-MDP M_j^f ;

Compute the feasibility $w(\xi_j)$ for each trajectory $\xi_j \in \Xi_j$ as in Eqn. (3);

Assign the feasiblity $w(\xi_j)$ to each state transitions in ξ_j ;

end

$$\text{Compute } p_w((s_t^d, s_{t+1}^d)) \leftarrow \frac{w((s_t^d, s_{t+1}^d))}{\sum_{\left(s_{t'}^d, s_{t'+1}^d\right) \in T} \frac{w((s_{t'}^d, s_{t+1}^d))}{w((s_{t'}^d, s_{t'+1}^d))}}$$

while not converging do

Sample a batch of state transitions from p_w ;

Train π with the sampled batch of state transitions by an state-based imitation learning algorithm: state-based GAIL

$$\min_{\theta_{\pi}} \max_{\theta_{d}} \mathbb{E}_{(s_{t}, s_{t+1}) \sim \pi} (\log(D(s_{t}, s_{t+1}))) + \mathbb{E}_{(s_{t}, s_{t+1}) \sim p_{w}} (1 - \log(D(s_{t}, s_{t+1}))), \tag{7}$$

where θ_{π} is the parameter of π and θ_d is the parameter for the discriminator D in state-based GAIL;

end

Output: Learned optimal policy π^* for \mathcal{M}^i .

Currently, we consider the state transitions in each trajectory share the same feasibility but do not consider the case where parts of the trajectories are more feasible than some other parts. This is because the state-based GAIL in our algorithm as well as many standard imitation learning algorithms rely on learning from the full trajectory from the start to the end state. If a segment of the trajectory is far from feasible or harmful, then the remaining part is also not going to be useful for our algorithm. Therefore, we only learn from trajectories that are helpful in all parts.

C Experiment Details

In this section, we discuss experimental details for all of our environments. In addition, we discuss the choice of our discount factor γ^f in C.3

C.1 Mujoco

Table 1: The composition of demonstrations for each environment

Environment	Number of Demonstrations				
2111011110111	i	ii	iii	iv	
Swimmer	50	50	500	500	
Walker2d (first setting)	50	50	500	500	
Walker2d (second setting)	500	500	10	10	
HalfCheetah (first setting)	25	25	500	500	
HalfCheetah (second setting)	500	500	25	25	
Hopper	50	50	500	500	

Composition of Demonstrations. The composition of demonstrations in each environment is shown in Table 1. We design the composition of demonstrations to ensure that directly performing imitation learning on all the demonstrations cannot learn an optimal policy as otherwise we do not need to consider the problem of learning from demonstrations from agents with different dynamics.

Implementation Details. To implement our algorithm, we use TRPO [31] as the RL algorithm to learn the optimal policy from f-MDP, and we use the GAIL from Observation algorithm [16] as our imitation learning technique to learn the optimal policy from the reweighted demonstrations. For each demonstrations, we create a separate f-MDP for its demonstrations and train an optimal policy for the f-MDP to generate the feasibility for its demonstrations.

Compared to the final imitation learning algorithm, which requires about 7×10^7 interactive steps with the environment to converge, learning the optimal policy from f-MDP only needs about 5×10^5 time steps, which is significantly smaller. This indicates that the proposed feasibility learning is efficient even with an additional RL learning process.

C.2 Simulated and Real Robot Arm

Reward Function. In the main text, we evaluate our policy for the two robot arm environments using the expected return. We introduce the details of the reward function here. The exact formula of the reward function is r=-s-10000h+5000g, where r is the reward, s is the number of steps, $h \in \{0,1\}$ represents whether the robot hits any object in the environment, and $g \in \{0,1\}$ represents whether the robot achieves the goal.

Composition of Demonstrations. The demonstration set composed of 5 trajectories from the 3 DoF Panda robot arm with disabled joints and 43 trajectories from the 7-DoF Panda arm for both the simulated and the real arm environments.

Implementation Details. To implement our algorithm, we use TRPO [31] as the RL algorithm to learn the optimal policy from f-MDP. To learn the optimal policy from the reweighted demonstrations, we learn a beta-VAE [32] to imitate the state transition sampled from the reweighted distribution of state transitions p_w . After learning the state transitions, we recover the joint actions from the changes of the end-effector's position using inverse kinematics.

Compared to learning the beta-VAE model, which requires about 2.56×10^4 interactive steps with the environment to converge, learning the optimal policy from f-MDP only needs about 5.12×10^3 time steps, which is negligible with respect to the imitation learning algorithm. This indicates that the proposed feasibility learning is efficient even with an additional RL learning process.

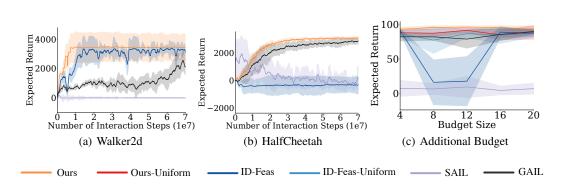


Figure 6: (a-b) show the expected return for the second setting for the Walker2d and the HalfCheetah environments respectively. (c) The expected return with varying budget of additional demonstrations in Swimmer.

C.3 Choice of Discount Factor

The choice of the discount factor γ^f depends on how fast the compounding error increases with respect to the number of steps in the environment. Faster increasing compounding errors need smaller γ^f . In practice, this depends on the scale of the 'movement' of the agent at each step. For example, for a robot arm, if each joint can only move at a small angle at each step, we can set γ^f to be larger or if each joint can move at a larger angle at each step, γ^f should be set smaller. In our experiments, we fix the $\gamma^f=0.9$ and empirically it works well for all the environments.

D Experimental Results

In this section, we provide a number of additional experiments and results including additional results in the Mujoco environment and the simulated robot (D.1, D.3), results demonstrating the effects of varying the number of demonstrations (D.4), discussion and additional results with varying choices of distance metrics (D.5), comparison with a mapping-based method (D.6), and comparison of the expected return of our method with the expected return of the originally collected demonstrations (D.7).

D.1 More Results for Mujoco Environments

We show the results of the second setting for the Walker2d and the HalfCheetah in Fig. 6. We observe that our proposed feasibility achieves the best performance among all the methods. The highest p-value comparing our method to baselines is 0.297 with ID-Feas for the Walker2d environment, and 0.0037 with ID-Feas for the HalfCheetah environment (statistically significant) respectively.

D.2 Performance with a Budget of Additional Demonstrations.

We now consider a setting, where we start with a limited set of demonstrations, but acquire more demonstrations under a limited budget. Our feasibility metric can assess how likely it is for a demonstrator to produce feasible demonstrations, and hence can help us select which demonstrator to query for more demonstrations. We start with one demonstration from each demonstrator in the Swimmer environment and evaluate the performance as we add demonstrations. For our method and ID-Feas, we can acquire demonstrations proportional to the computed feasibility score. We compare the expected return with demonstrations selected based on feasibility (Ours, ID-Feas) to the expected return with demonstrations uniformly acquired from each demonstrator (Ours-Uniform, ID-Feas-Uniform). We further compare with SAIL and GAIL, where no feasibility is defined and we uniformly acquire demonstrations. As shown in Fig. 6(c), Ours outperforms ID-Feas, which demonstrates that the proposed feasibility can better reflect how likely each demonstrator produces feasible

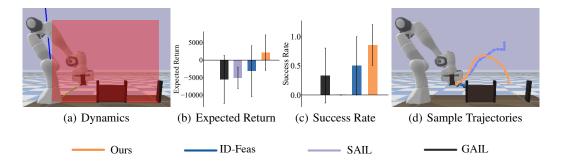


Figure 7: (a) The illustration of different dynamics in the simulated robot arm environment. The 7 DoF robot arm can move in the whole 3D space while the 3 DoF arm can only move in the red plane; (b-c) The bar plot for the expected return and the success rate for the simulated robot arm environment; (d) Sampled trajectories for different methods in the simulated robot arm environment.

demonstrations and acquire more demonstrations from helpful demonstrators. Ours outperforms all the other methods including Ours-Uniform (although not with statistical significance), which indicates that the demonstrations acquired based on the feasibility gain more useful information.

D.3 Results on the Simulated Robot

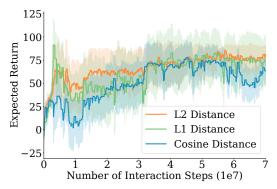
As shown in Fig. 7, We observe that the proposed approach outperforms the baseline methods both in terms of the expected return and the success rate. The highest p-value for the expected return is 7.252×10^{-9} and for the success rate is 1.047×10^{-8} (both with ID-Feas), which are both statistically significant. The sampled trajectory show that the proposed approach achieves an efficient trajectory successfully moving the book to the empty area of the shelf.

Table 2: The performance of varying percentage of demonstrations used with respect to the original setting.

Method	sthod Swimmer		Walker2d			HalfCheetah			Hopper			
	20%	50%	100%	20%	50%	100%	40%	60%	100%	20%	50%	100%
GAIL	40.0±34.3	37.9±35.2	30.9±23.5	276.9±18.3	313.1±63.0	261.1±5.6	2464.9±460.2	2597.2±399.0	2443.6±440.7	2798.3±351.1	2996.6±623.2	3009.6±362.4
SAIL	-7.7 ± 19.7	-5.5±24.6	-3.0 ± 28.4	8.5 ± 17.5	-5.3 ± 56.4	19.0 ± 30.1	-556.7±365.8	-503.3 ± 299.3	-603.0±389.6	-252.6±432.6	-1622.2±1780.1	-7.00 ± 11.4
RAL	23.1 ± 33.9	30.6 ± 28.1	48.2 ± 39.0	244.3±27.4	261.0 ± 46.7	227.0 ± 24.2	2604.4±423.1	2515.3 ± 311.0	2594.4±508.7	2040.3±408.3	2108.9±611.9	1916.1±750.4
Ours	68.2 ±24.7	76.4 ±22.1	$\textbf{74.3} {\pm} 20.1$	3137.6 ±50.2	3127.0 ±30.7	3144.3 ±23.5	2716.7±301.6	2812.4 ±261.2	2830.6 ±292.6	3273.3±180.2	3351.0 ±146.3	3329.6 ±115.2

D.4 Varying the Number of All Demonstrations

We vary the number of demonstrations from all demonstrators. We conduct experiments on the first setting of the all the Mujoco environments. For the Swimmer and Walker2d environments, we test with 20% and 50% of the original demonstrations. For the HalfCheetah environment, we test with 40% and 60% of the original demonstrations since we have much less demonstrations (25 vs 50) from the demonstrators similar to the imitator. As shown in Table 2, we observe that our approach outperforms all the other methods when having access to different number of demonstrations.



Method	Expected Return
GAIL [4]	31.20±22.25
SAIL [5]	0.56 ± 4.27
DCC [26]	5.32 ± 3.43
ID-Feas [7]	48.96 ± 38.50
Ours	74.89 ± 19.68

Figure 8: The expected return with respect to the Table 3: The expected return of the learned polmetrics.

number of steps with different choices of distance icy in the Swimmer environment (with standard deviation).

Discussion of the Choice of Distance Functions

We use L2 distance between states in the reward function in Eqn. (2) and (3) in the main text and in all the experiments. This is because in all our environments, the L2 distance can accurately measure the distance between states. However, this does not mean that the distance metric in our reward of f-MDP is restricted to the L2 distance. We can change the distance depending on the specific state space. For example, for a state space with unit vectors, we can use the cosine distance as the distance metric.

In Fig. 8, we show the expected return of our method by using different distances in the Swimmer environment. We use L1 distance and Cosine distance (the cosine of the angles between two state vectors) as examples. We observe that L1 distance, which is another distance derived by norm, performs close to L2 distance, but Cosine distance performs worse than L2 distance because Cosine distance only cares about the distance on angle but ignores the scale of the vectors, while in Swimmer, the scale of the states is also important. The results show that the choice of this distance function is flexible and depends on the specific choice of the state space in our problem.

D.6 Comparison with Mapping-Based Methods

Mapping-based methods translate the demonstrations across different environments by learning state mappings and action mappings [26], which can be used in our problem setting by mapping the source demonstrations to the target environment. However, our problem setting does not ensure that there exist a mapping between the demonstrators and the imitator, which violates the assumption of the mapping-based methods. We thus do not include any mapping-based methods in our experiments in the main body of our work. However, here as an additional experiment, we compare our method with the state-of-the-art mapping-based method, DCC [26].

DCC requires random trajectories from both the demonstrators and the imitator to learn a mapping, but we do not have access to the demonstrators' environment and only have access to demonstrators' demonstrations. So we use the demonstrations and the imitator's random trajectories as the input to DCC. As shown in Table 3, the performance of DCC is much worse than our method and even worse than GAIL. This is because DCC itself is a good mapping-based method but mapping-based methods are not quite suitable for our problem setting. In fact, there should not exist a mapping between demonstrations and the imitator's random trajectories. Building such a mapping causes severe mismatch between states and actions of different environments and makes the translated demonstrations distort the original demonstrations.

Comparison with the Collected Demonstrations

We compare the expected return of our approach with the demonstrations in Table 4. We observe that in the first setting of Walker2d environment, Hopper environment, the simulated robot and the real

Table 4: The average expected return of demonstrations in different environments and the expected return of our method.

	Swimme	Walk	Walker2d		HalfCheetah		Simulated	Real
	Swimme	First	Second	First	Second	Hopper	Robot	Robot
Demonstrations	106±3	3098±118	3720±336	3229±170	3337±67	3460±87	1823±110	2531 ± 362
Ours	75±20	3147 ± 10	3424 ± 645	$2832{\pm}291$	$3142{\pm}89$	$3330 {\pm} 115$	2127 ± 5053	2746 ± 2712

robot environments, our approach performs comparably to the expected return of demonstrations, which are optimal demonstrations for different demonstrators. In the Swimmer, the second setting of Walker2d and the HalfCheetah environments, the performance is worse than the demonstrations. This is because only a few demonstrations are feasible for the imitator and that may not be enough to learn an optimal policy. However, the margin between our approach and the demonstrations is still not large. The results show that the proposed feasibility can select useful demonstrations for the imitator to imitate.